

Concentration of the number of solutions of random planted CSPs and Goldreich's one-way function candidates

Emmanuel Abbe^{*}

Katherine Edwards[†]

May 1, 2015

Abstract

This paper shows that the logarithm of the number of solutions of a random planted k -SAT formula concentrates around a deterministic n -independent threshold. Specifically, if $F_k^*(\alpha, n)$ is a random k -SAT formula on n variables, with clause density α and with a uniformly drawn planted solution, there exists a function $\phi_k(\cdot)$ such that, besides for some α in a set of Lebesgue measure zero, we have $\frac{1}{n} \log Z(F_k^*(\alpha, n)) \rightarrow \phi_k(\alpha)$ in probability, where $Z(F)$ is the number of solutions of the formula F . This settles a problem left open in Abbe-Montanari RANDOM 2013, where the concentration is obtained only for the expected logarithm over the clause distribution. The result is also extended to a more general class of random planted CSPs; in particular, it is shown that the number of pre-images for the Goldreich one-way function model concentrates for some choices of the predicates.

^{*}Program in Applied and Computational Mathematics, and EE Department, Princeton University, Princeton, NJ. Email: eabb@princeton.edu.

[†]Department of Computer Science, Princeton University, Princeton, NJ. Email: ke@princeton.edu.

1 Introduction

This paper investigates concentration phenomena for the number of solutions in random planted random constraint satisfaction problems (CSPs) and the Goldreich one-way function candidate.

A large body of works have studied phase transition phenomena for satisfiability in random CSPs. For uniform¹ models, the probability of being satisfiable often tends to a step function as n tends to infinity, jumping from 1 to 0 when the constraint density crosses a critical threshold. For random k -XORSAT the existence of such a critical threshold is proved [21, 23, 25, 46]. For random 2-SAT, the threshold is proved in [17, 22, 33]. For random k -SAT, $k \geq 3$, the existence of an n -dependent threshold is proved in [31], and the satisfiability threshold conjecture states that this threshold is n -independent for all k . Recently, the conjecture was settled for k large enough [24], while upper and lower bounds are known to match up to a term that is of relative order $k 2^{-k}$ as k increases [11, 18]. Moreover, phase transition phenomena were also studied for a broad family of other CSPs, see for example [10, 11, 44] and references therein.

The counting problem for random formulas has also received attention recently. In [4], a concentration result is obtained for the number of solutions: at a fixed clause density α , the number of solutions of a random 2-SAT formula concentrates in the logarithmic scale to a deterministic n -independent threshold for almost every α . This result is extended for $k \geq 3$ for all clause densities having an UNSAT probability decaying fast enough (with a mild logarithmic decay being enough), which is conjectured to take place up to the SAT threshold. This result is obtained in two parts. First, as was shown earlier in [2, 6], the property that a random k -SAT formula has a number of solution bounded by 2^{n^ϕ} , for a fixed ϕ , has a phase transition with an n -dependent threshold, proved à la Friedgut. This is then turned into a concentration result for the number of solutions in [4] by showing that the limit for $\frac{1}{n} \mathbb{E} \log(1 + Z(F))$ exists, where $Z(F)$ is the number of solutions of the random formula. Observe that this gives an n -independent threshold for the concentration. The key tool in establishing this limit is the interpolation method, first introduced in [37] for the Sherrington-Kirkpatrick model, and subsequently generalized and extended in [4, 15, 28, 29, 45]. Note however that the use of “1+” in the logarithm above (to obtain a well defined quantity) is responsible for the difficulty in obtaining the concentration for all clause densities when $k \geq 3$.

In this paper, we consider CSPs that have a planted solution and study the counting problem for random ensembles. Planted CSPs are a rich ground for studying combinatorial optimization problems motivated by ‘real-world’ applications, such as in coding theory, community detection, or cryptography, where a solution typically does exist but where the problem is to identify how many other solutions are there, or how hard is it to recover the planted solution. Planted ensembles were investigated in [5, 8, 9, 14, 38, 39], and at high density in [12, 19, 27], and relationships between planted random CSPs and their non-planted counterparts in the satisfiable phase were studied in [5, 42, 48].

It was shown recently in [3] that for a broad class of random planted CSPs, the logarithm of the number of solution concentrates to an n -independent deterministic threshold for almost every clause density. In particular, this covers k -SAT for all k ’s. Hence, the planting allows to circumvent the issues of establishing the limit of the log-partition function, since the latter is well defined due to the planting (no need for the “1+” term discussed above). However, the planting also introduces asymmetry in the model, which lead [3] to a weaker concentration result: the concentration is obtained with respect to the graph ensemble but is taken *in expectation* over the clause distribution.

¹The model may have a fixed but uniform number of constraints, or a Binomial or equivalent form.

Let us explain this nuance more precisely for random k -SAT. A random planted formula is defined in this case by drawing first a random uniform solution x^0 , and independently, a random 3-hypergraph $G = ([n], E)$ at a fixed edge density. The random clauses are then defined for each edge $e \in E(G)$ by drawing a negation pattern s_e uniformly at random within the set of negations patterns that preserve x^0 as a planted solution. Specifically, the clause for edge e is defined by $y[e] \neq s_e$ (where $y[e]$ is an assignment of literals to the variables associated with e). Note that this is indeed equivalent to requiring that the OR of the variables in $y[e]$ negated with the pattern x_e is 1. Consider now

$$\phi_n := \frac{1}{n} \log(Z(F^{(0)})),$$

where $F^{(0)}$ is the random planted formula. In [1], it is shown that $\mathbb{E}_s \phi_n$, the expectation of ϕ_n taken over the variables $s = \{s_e\}_{e \in E(G)}$, concentrates in probability (with respect to the drawing of G) to a deterministic n -independent value.² It was left open to obtain concentration with respect to the drawing of s as well. In particular, the martingale argument used in [1] fails in this case, since the fluctuations are not bounded, and the application of Friedgut's theorem is mitigated by the lack of symmetry caused by the planting.

We resolve in this paper the above problem left open in [1] and show that for almost every α , there exists an n -independent value $\phi(\alpha)$ such that

$$\frac{1}{n} \log(Z(F^{(0)})) \rightarrow \phi(\alpha) \quad \text{in probability,}$$

closing the concentration problem. The main tool is based on Bourgain's result from the appendix of [31]. The result is then generalized to a broad class of planted CSPs, and a new application to Goldreich's one way function [36] is investigated.

The Goldreich one-way function candidate is defined from a k -hypergraph G on n vertices and m hyperedges and a *fixed* predicate function $\chi : \{0, 1\}^k \rightarrow \{0, 1\}$. The function takes an input $x \in \{0, 1\}^n$ and, evaluating χ at each of the m k -tuples selected by the hyperedges of G , produces an output in $\{0, 1\}^m$. In [36], G is proposed to be drawn at random with an edge density m/n , and the choice of predicates is further discussed in [16]. Note that both $m = \omega(n)$ and $m = \Theta(n)$ are potential candidates [16, 36]. Defining the rate of the one-way function by m/n , it is interesting to understand for what rates (in addition to what predicates) is the function possibly one-way, in particular, for the case of $m/n = \alpha$ constant. A natural approach would be to relate this question to the structure of the solution space of the underlying CSP, starting with its size, and hypothetically with the condensation [20, 41] and freezing of the solution clusters [47] phenomena.³ In particular, the function $\phi(\cdot)$ is expected to have a kink at the condensation threshold for various CSPs [43], and this may indicate a behavioral changes for the hardness of the one-way function. In this paper, we investigate the most basic question towards such considerations: does the function ϕ even exist? Namely, does the normalized logarithm of the number of pre-images concentrates for some/all predicates?

We answer this question by the affirmative for a certain class of predicates. Interestingly, it is not obvious that this class of predicates overlaps with the class of predicates that precludes the non-hardness conditions introduced in [16] for large clause densities. We hence leave an open problem: can one obtain concentration and hardness at the same time, or is hardness related to the non-concentration? We believe that the former is true and that our proof

²Note that the variables $s = \{s_e\}_{e \in E(G)}$ depend on the planted assignment. For a deterministic kernel Q , this means in expectation over the planted assignment.

³In the non-planted models, these phenomena have been associated with computational barriers for satisfiability.

technique stumbles on technicalities, but we cannot resolve this argument.

2 Models

2.1 CSPs arising from satisfiability problems

We first describe a class of constraint satisfaction problems. Let $V = \{v_1, \dots, v_n\}$ be a set of Boolean variables, and fix an integer $k \geq 2$. An instance F of a CSP consists of a k -uniform multi-hypergraph (V, E) (that is, all edges have cardinality k and we allow parallel edges), and a family of *clause functions* $\chi_e : \{0, 1\}^k \rightarrow \{0, 1\}$ for each $e \in E$. A k -*clause* comprises an edge e and its corresponding function χ_e . We'll sometimes call F a formula. The form of the clause function depends on the type of satisfiability problem we are interested in (for the moment, SAT, NAESAT or XORSAT). Let $y[V]$ denote an assignment y_1, \dots, y_n of Boolean values to the variables in V , and $y[e]$ its restriction to the k variables in e . By $\chi_e(y[e])$ we mean the result of evaluating χ_e on the k values in some *fixed* order (for the moment, the actual choice of ordering of variables in edges isn't important but it will be when we consider certain planted models in Section 4.2.) This model naturally captures familiar satisfiability problems:

- in k -SAT, we have $\chi_e(y[e]) = 1 \iff y[e] \neq x_e$ where $x_e \in \{0, 1\}^k$ represents a particular forbidden pattern,
- in k -NAESAT, we have $\chi_e(y[e]) = 1 \iff y[e] \notin \{x_e, \bar{x}_e\}$ where \bar{x}_e is the result of flipping each component of x_e ,
- in k -XORSAT, we have $\chi_e(y[e]) = 1 \iff \oplus_i y_i = x_e$ where now $x_e \in \{0, 1\}$.

An assignment which satisfies all clauses in F is called a satisfying assignment (or solution) for F . Let $C_k(n)$ be the set of all possible k -clauses on V and write $N = |C_k(n)|$; in k -SAT for example we have $N = \binom{n}{k} 2^k$. We use the *binomial model* for a random CSP with clause density $\alpha \in [0, N/n]$, and draw a random formula G as follows: ⁴

- include in G each clause in $C_k(n)$ with probability $p = \alpha n/N$. (1)

Let $G(n, \alpha)$ denote a formula obtained by this process. The formula $G(n, \alpha)$ can be viewed as a random element of $\{0, 1\}^N$, drawn according to the product measure μ_p . ⁵ That is, for each $x \in \{0, 1\}^N$ we have $\mu_p(x) := \mathbb{P}[G(n, \alpha) = x] = p^{|x|} (1-p)^{N-|x|}$ (where $|x|$ denotes the number of nonzero components). Now consider a procedure to sample a *planted* CSP F .

- Sample $v^0 \in \{0, 1\}^n$ uniformly at random. (2)
- Then include in F each k -clause which is satisfied by v^0 independently with probability $p = \alpha n/N$.

We use the notation $F(n, \alpha)$ to denote a formula obtained by this process. By construction such a formula is always satisfied by the assignment $v_i = v_i^0$; the vector v^0 is known as the planted solution. Let $Z(F)$ denote the cardinality of the set of assignments to v_1, \dots, v_n which satisfy F . Notice that we always have $Z(F(n, \alpha)) \geq 1$ by construction. Again we

⁴One could also consider the *uniform model*, wherein $G(n, \alpha)$ is chosen uniformly from those vectors $x \in \{0, 1\}^N$ with $|x| = \alpha n$, where $\alpha n/N = p$. The models are essentially equivalent, and we mostly focus here on the binomial model.

⁵To see the correspondance, identify each of the N components with a clause and set it to 1 if and only if the clause is present in the formula.

view $F(n, \alpha)$ as an element of $\{0, 1\}^N$ but observe that the distribution in this case satisfies $\mu_p(F) = \mathbb{P}[F(n, \alpha) = F] = \frac{Z(F)}{2^n} p^{|F|} (1-p)^{\binom{n}{k}(2^k-1)-|F|}$ so μ_p is not a product measure here.

2.2 CSPs arising from Goldreich's one-way function candidate

The CSPs introduced in the previous section have clause functions taking a few specific forms. In these examples a satisfying variable assignment $y[V]$ satisfies $\chi_e(y[e]) = 1, \forall e$, and the clause functions on individual edges are independent of one another. Our concentration results can be extended to a related class of CSPs which are related to Goldreich's proposed one-way function [34]. The idea is that we can consider CSPs with arbitrary clause functions if the clauses on different edges are related in a specific way.

In [34] Goldreich proposed a candidate one-way function family which exploits the difficulty of recovering a solution to a form of planted CSP.⁶ Goldreich's original proposition was that the following function f is one-way. As always we work with the variable set $V = \{v_1, \dots, v_n\}$.

- Select a predicate $\chi : \{0, 1\}^k \rightarrow \{0, 1\}$ uniformly at random from the set of all such Boolean functions.
- Draw a sparse Erdős-Rényi k -uniform multi-hypergraph (V, E) with m edges e_1, \dots, e_m .
- $f : \{0, 1\}^n \rightarrow \{0, 1\}^m$ is the function with $f(x)_i = \chi(x[e_i])$, i.e. the i th output bit is the result of evaluating ϕ on the k (ordered) values assigned to the edge e_i .

More precisely, Goldreich conjectured that f is one-way in the setting where $k = O(\log n)$ and $m = n$, and the graph is a sufficiently good expander, for most choices of the predicate χ which is randomly selected and hard-wired into f .

With this in mind we can define a class of planted CSPs generated by the following procedure.

- Select a predicate $\chi : \{0, 1\}^k \rightarrow \{0, 1\}$.
- Sample $v^0 \in \{0, 1\}^n$ uniformly at random.
- Then include in F each k -clause of the form e with $\chi(y[e]) = \chi(v^0[e])$, with probability $p = \alpha n/N$.

(3)

Here the edges are ordered subsets of V , and so $N = \binom{n}{k} k!$.

3 Overview of results

Recall that for is a CSP formula F (planted or not) we denote by $Z(F)$ the number of satisfying assignments for F . If $\phi \in [0, 1]$ we write $Q_n(\alpha, \phi) := \mathbb{P}[Z(F(n, \alpha)) < 2^{n\phi}]$.

3.1 Concentration of the number of solutions of planted satisfiability CSPs

Our main result is the following theorem, which states that for fixed $\alpha \geq 0$ the logarithm of the number of solutions of a random planted formula concentrates, closing the problem left open in [4]. Note that this clears the concentration problem in its most general form: the exponent of the number of solutions of a random planted SAT formula can be asymptotically predicted with an n -independent value and for any $k \geq 2$ (small or large). The only part that could be further generalized is the fact that the result does not hold for a countable set of

⁶One-way functions are important objects in cryptography and complexity theory. Intuitively these are functions that are computationally easy to evaluate, but hard to invert. For a thorough discussion see [35], [34].

“bad” α ’s, but it is unclear whether this is a technicality or not. The formal result reads as follows.

Theorem 1. *For every $k \geq 2$, there exist a countable set \mathcal{D} and a function $\phi_s : [0, \alpha^*] \rightarrow [0, 1]$ such that for every $\alpha \notin \mathcal{D}$ and every $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} Q_n(\alpha, \phi_s(\alpha) - \epsilon) = 0$$

$$\lim_{n \rightarrow \infty} Q_n(\alpha, \phi_s(\alpha) + \epsilon) = 1$$

In [3], it was shown that this quantity concentrates when the expectation is taken over the clause distribution. We use this result in the proof of Theorem 1.

Theorem 2. [3] *For every $k \geq 2$, for every $\alpha \in [0, \alpha^*]$ the sequence*

$$\psi_n(\alpha) := \frac{1}{n} \mathbb{E}[\log Z(F(n, \alpha))]$$

converges almost surely to a limit $\phi_s(\alpha)$.

As an intermediate step toward Theorem 1, we will prove that for fixed $\phi \in [0, 1]$ there is a sharp threshold density for the property of having fewer than $2^{n\phi}$ solutions (we define these terms in Section 4.1). First, in Section 4.2 we prove the following n -dependent sharp threshold.

Lemma 3. *For every $k \geq 2$ and for every $\phi \in [0, 1)$ there exists a sequence $\{\alpha_n(\phi)\}_{n \in \mathbb{Z}_{>0}}$ such that for every $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} Q_n(\alpha_n(\phi) - \epsilon, \phi) = 0$$

$$\lim_{n \rightarrow \infty} Q_n(\alpha_n(\phi) + \epsilon, \phi) = 1.$$

In fact, we prove Lemma 3 for a larger class of planted CSPs, namely those which arise from Goldreich’s one-way function candidate [34]. This allows us to deduce the analogous statement of Theorem 1 for certain instances of these CSPs, as well as an n -dependent version of it in general.

In Section A we combine Lemma 3 with Theorem 2 using a technique from [4] to show that the sequence $\alpha_n(\phi)$ converges.

Theorem 4. *For every $k \geq 2$, there exist a countable set \mathcal{C} and a function $\phi_s : [0, \alpha^*) \rightarrow [0, 1]$ such that for each $\phi \in \phi_s([0, \infty))$ there exists $\alpha_s(\phi)$ such that for each $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} Q_n(\alpha_s(\phi) - \epsilon, \phi) = 0$$

and

$$\lim_{n \rightarrow \infty} Q_n(\alpha_s(\phi) + \epsilon, \phi) = 1$$

We deduce Theorem 1 from Theorem 4 in Section A.

3.2 Concentration of the number of solutions of CSPs from Goldreich’s one-way function candidates

We now present concentration results for the number of solutions of the CSPs arising from Goldreich’s one-way function candidates described in Section 2.2.

If one considers the logarithm of the number of solutions of the one-way function candidate determined by a random graph G , a predicate χ and a uniform input, and takes the average

over the input distribution, it is possible to obtain the following concentration result. Note that this gives a stronger concentration notion, i.e., almost sure and for every α , and imposes no restriction on the choice of χ . However, it provides a n -dependent threshold and requires averaging over the input distribution.

Lemma 5. *Let $F(n, \alpha)$ be a formula drawn as in (3). Then for every $k \geq 2$, there exist a function $\phi_s^n : [0, \alpha^*] \rightarrow [0, 1]$, namely $\phi_s^n = \mathbb{E}_{G, v^0} \log Z(F(n, \alpha))$, such that for every $\alpha > 0$ and every $\epsilon > 0$, the following holds almost surely*

$$\lim_{n \rightarrow \infty} (\mathbb{E}_{v^0} \log Z(F(n, \alpha)) - \mathbb{E}_{G, v^0} \log Z(F(n, \alpha))) = 0.$$

The proof is found in Appendix Section B.

We can dispose of the dependence of ϕ_s on n and on the averaging of the input in the previous theorem for certain choices of χ . We simply need to remark that Theorem 2 was in fact shown in [4] to hold for planted formulas $F(n, \alpha)$ which satisfy a certain convexity hypothesis (let's call it H for now), then the proof of the following theorem follows that of Theorem 1 exactly as in Section A. To this end, in the Appendix Section C we prove the analogue of Lemma 3 for these CSPs. This allows us to deduce the analogous statement of Theorem 1 for certain instances of these CSPs.

We say that a predicate $\chi : \{0, 1\}^k \rightarrow \{0, 1\}$ is *balanced* if it evaluates to 1 on exactly half of the inputs and we say χ is *antisymmetric* if $\chi(x) = 1 - \chi(\bar{x})$ for some $x \in \{0, 1\}^k$.

Theorem 6. *Let $F(n, \alpha)$ be a formula drawn as in 3, with a predicate χ which is antisymmetric and satisfies Hypothesis H. Then for every $k \geq 2$, there exist a countable set \mathcal{D} and a function $\phi_s : [0, \alpha^*] \rightarrow [0, 1]$ such that for every $\alpha \notin \mathcal{D}$ and every $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} Q_n(\alpha, \phi_s(\alpha) - \epsilon) = 0$$

$$\lim_{n \rightarrow \infty} Q_n(\alpha, \phi_s(\alpha) + \epsilon) = 1$$

The hypothesis H , stated in terms of χ is as follows.

Definition 1. *Let $M_1(\{0, 1\}^k)$ denote the space of probability measures on $\{0, 1\}^k$. Let $\ell \geq 1$. Define $\Gamma : M_1(\{0, 1\}^k) \rightarrow \mathbb{R}$ by*

$$\nu \mapsto \Gamma_\ell(\nu) = \frac{1}{2} \sum_{\substack{u^{(1)}, \dots, u^{(\ell)} \in \{0, 1\}^k \\ \chi(u^{(1)}) = \dots = \chi(u^{(\ell)})}} \prod_{i=1}^k \nu(u_i^{(1)}, \dots, u_i^{(\ell)})$$

Hypothesis H. *For each $\ell \geq 1$, the operator Γ is convex in ν .*

Bogdanov and Qiao showed in [16] that for many choices of χ , Goldreich's function can be inverted with high probability when m is larger than n by a sufficiently large constant factor. In particular any χ which is not balanced or whose output correlates with one or two bits of the input is a bad choice when $m = Dn$ for sufficiently large constant D . Their result suggests that if we want the resulting function to be one-way then we may want χ to be balanced and not correlated with any bit or pair of bits of the input, but it is unclear whether these would be necessary in the regime $m = n$, the one Goldreich originally suggested.

Strictly speaking, the restriction to antisymmetric χ in Theorem 6 does not seem necessary. It is a technical condition which arises in the proof of Lemma 3. We have verified using a computer search that when $k \leq 5$ no antisymmetric function satisfies the balance properties along with Hypothesis H but it remains unclear to us whether such a function can exist in general.

4 An overview of the proofs

The main element in our proofs is Lemma 3, whose proof we give in this section. From there, obtaining Theorem 1 and its analogues is a straightforward argument given in the Appendix Section A.

4.1 Sharp thresholds and Bourgain's theorem

Before proceeding to the proof of Lemma 3, we briefly give a bit of background material on sharp thresholds. A subset $\mathcal{A}_n \subseteq \{0, 1\}^N$ is called a *property*, and we say it is *nontrivial* if $\mathcal{A}_n \subset \{0, 1\}^N$. Property \mathcal{A}_n is *monotone increasing* (or simply *monotone*) if for every $x \in \mathcal{A}$ and $x \subseteq y$ we have $y \in \mathcal{A}$. (Containment of formulas is defined in the natural way, namely $x \subseteq y$ iff every nonzero component of x is also nonzero in y .) We may drop the subscript n when it is unambiguous or unnecessary. A property is *symmetric* if there is a transitive permutation group under which it is invariant. For example, in (unplanted) SAT, the property of being unsatisfiable is monotone and symmetric.

In this section and the next it is convenient to make a slight abuse of notation, and write $F(n, p)$ in place of $F(n, \alpha)$ to stress that clauses are included in $F(n, \alpha)$ according to binomial ($p = \alpha n/N$) distribution. For a monotone property $\mathcal{A}_n \subset \{0, 1\}^N$, write $\mu_p(\mathcal{A}_n) = \sum_{x \in \mathcal{A}_n} \mu_p(x) = \mathbb{P}[F(n, p) \in \mathcal{A}_n]$. It's not difficult to show that if \mathcal{A} is a nontrivial property then $\mu_p(\mathcal{A})$ is a strictly increasing and continuous function of p . For $\gamma \in (0, 1)$, let $p_n(\gamma)$ be the value which (uniquely) satisfies $\mu_{p_n(\gamma)}(\mathcal{A}) = \gamma$. We say that \hat{p}_n is a threshold probability if

$$\lim_{n \rightarrow \infty} \mathbb{P}[F(n, p) \in \mathcal{A}] = \begin{cases} 1 & \text{if } p_n \gg \hat{p}_n \\ 0 & \text{if } p_n \ll \hat{p}_n \end{cases}$$

where the notation $p_n \gg \hat{p}_n$ indicates that $\frac{\hat{p}_n}{p_n} \rightarrow 0$ as n diverges.

We make a distinction between properties exhibiting a very rapid transition versus those with a more gradual one. Formally, we say that \mathcal{A} has a sharp threshold if for every $\gamma \in (0, 1)$ there exists $p_\gamma = p_\gamma(n)$ such that $\mathbb{P}[F(n, p_\gamma) \in \mathcal{A}] = \gamma$, and such that for every $\delta > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}[F(n, p) \in \mathcal{A}] = \begin{cases} 1 & : \text{ if } p(n) \geq (1 + \delta)p_\gamma(n) \\ 0 & : \text{ if } p(n) \leq (1 - \delta)p_\gamma(n) \end{cases}$$

Equivalently, for $\tau \in (0, 1)$ define p_0, p_1, p_c such that $\mu(p_0) = \tau$, $\mu(p_1) = 1 - \tau$ and $\mu(p_c) = \frac{1}{2}$. The property \mathcal{A} has a *sharp* threshold if the ratio $\frac{p_1 - p_0}{p_c}$ tends to 0. The threshold is *coarse* if this ratio is bounded away from 0, i.e. if there exists some constant C such that for some $\gamma \in (0, 1)$ we have $p_\gamma \frac{d\mu_p(\mathcal{A})}{dp} \big|_{p=p_\gamma} < C$. Friedgut and Kalai (see [32]) showed that in this case it must be true that $p_\gamma = o(1)$.

A crucial contribution to the theory of sharp thresholds is due to Friedgut, in the form of a general existence theorem for sharp thresholds (see [30], [31]). Roughly, the theorem asserts that if a monotone symmetric property has a coarse threshold, then it can be approximated by the property of containing a small fixed subgraph. We omit the statement of Friedgut's theorem since it does not apply in our setting; introducing a planted solution does away with the symmetry in the properties we are interested in. Fortunately in the appendix to [31], Bourgain gave an analogue of Friedgut's result to nonsymmetric properties as follows. This is the theorem we will need to apply.

Theorem 7 (Bourgain [31]). (See also [40]) *Let $\mathcal{A}_n \subset \{0, 1\}^N$ be a monotone property, and $C > 0$ constant. Suppose μ_p is the product measure on $\{0, 1\}^N$, i.e. $\mu_p(x) = p^{|x|}(1-p)^{N-|x|}$ for every x . Assume that there exists $\gamma \in (0, 1)$ such that $\mu_{p_\gamma}(\mathcal{A}_n) = \gamma$ and $p_\gamma \frac{d\mu_p(\mathcal{A}_n)}{dp} \big|_{p=p_\gamma} < C$ and $p = o(1)$. Then there exists $\delta = \delta(C) > 0$ such that either*

1. $\mu_p(x \in \{0,1\}^n : x \text{ contains } x' \in \mathcal{A}_n \text{ of size } |x'| \leq 10C) > \delta$, or
2. there exists $x' \notin \mathcal{A}_n$ of size $|x'| \leq 10C$ such that the conditional probability satisfies

$$\mu_p(x \in \mathcal{A}_n | x' \subset x) > \gamma + \delta.$$

Friedgut's theorem and Theorem 7 provide a framework for finding sharp thresholds that has been widely exploited. These theorems typically allow one to prove the existence of a sharp threshold whose value depends on n , whereas in many cases the threshold is believed to converge. Friedgut's original application was to show that satisfiability for k -SAT has a sharp threshold. He also used the theorem to prove that in hypergraphs, the property of having a perfect matching, as well as 2-colourability have sharp thresholds. With Achlioptas in [7] they proved that k -colourability of graphs (for fixed k) has a sharp threshold. Krivelevich and Nachmias in [40] showed the same for list-colourability of bipartite graphs. Their proof uses a neat combinatorial trick (due to Alon) of combining Theorem 7 with a theorem of Erdős and Simonovits. We use a similar approach in the next section. A comprehensive survey of applications of Friedgut's theorem can be found in [30].

4.2 An n -dependent sharp threshold for planted CSPs

Here, we prove Lemma 3. For a fixed $\phi > 0$, we are interested in the property $\mathcal{A}_\phi (= \mathcal{A}_{\phi_n}) = \{F \in \{0,1\}^N; Z(F) < 2^{\phi n}\}$. Clearly, \mathcal{A}_ϕ is monotone increasing. We will show that it has a sharp (n -dependent) threshold. As before, let $\mathcal{A}_\phi = \{F \in \{0,1\}^N; Z(F) < 2^{\phi n}\}$, and now let $F = F(n, p)$ denote a CSP obtained as in (2). (We explain how the proof can be adjusted to handle $F(n, p)$ as in (3) in the Appendix Section C.) Lemma 3 can be restated as follows.

Lemma 8. *For a fixed k and $\phi > 0$, the property \mathcal{A}_ϕ has a sharp threshold.*

To prove Lemma 8 we will apply Bourgain's Theorem (Theorem 7). In the distribution of $F(n, p)$, we do not have the assumption on μ_p in the hypothesis of Theorem 7. To overcome this difficulty we need to consider fixed plantings, and observe that conditioning on the random choice of v^0 doesn't change the probability of the property \mathcal{A}_ϕ . By total probability,

$$\mathbb{P}[F(n, p) \in \mathcal{A}_\phi] = \sum_{v \in \{0,1\}^n} \mathbb{P}[F(n, p) \in \mathcal{A}_\phi | v^0 = v] \mathbb{P}[v^0 = v].$$

Further, for any $v \in \{0,1\}^n$, the conditional probability satisfies

$$\mathbb{P}[F(n, p) \in \mathcal{A}_\phi | v^0 = v] = \mathbb{P}[F(n, p) \in \mathcal{A}_\phi | v^0 = 0^n]$$

since the number of satisfying assignments is unchanged by swapping a variable with its negation.

Therefore, if we let $F^0(n, p)$ denote a formula obtained by independently including each k -clause which is satisfied by $v^0 = 0^n$ with probability p , we have $\mathbb{P}[F(n, p) \in \mathcal{A}_\phi] = \mathbb{P}[F^0(n, p) \in \mathcal{A}_\phi]$. So to prove Lemma 8 it is enough to show that \mathcal{A}_ϕ has a sharp threshold when 0^n is the planted solution. Now, the space we are working in is $\{0,1\}^{N'}$, where $N' = \binom{n}{k}(2^k - 1)$, and indeed $\mu_p(F) = p^{|F|}(1-p)^{N'-|F|}$. For the remainder of the proof, this will be the assumed setting.

We now proceed to prove the sharp threshold. The idea is to assume for a contradiction that \mathcal{A}_ϕ has a coarse threshold, and apply Bourgain's theorem. We closely follow arguments found in [40] and [6]. Roughly, Bourgain's theorem implies the existence of some fixed small formula x' whose appearance in a random formula increases the probability of having property

\mathcal{A}_ϕ by a positive amount. Note that \mathcal{A}_ϕ , while not symmetric, is invariant under relabelings of the variable set (i.e. automorphisms of $\{v_1, \dots, v_n, \neg v_1, \dots, \neg v_n\}$ which map $\{v_1, \dots, v_n\}$ to itself and $\neg v_i$ to the negation of the image of v_i , for each i). This property is sometimes called *permutation symmetry*. Thus, containing a random (relabelled) copy of x' has the same effect on the probability of having \mathcal{A}_ϕ . On the other hand, the assumption that the threshold is coarse implies that adding a large number of random clauses does not drastically change the probability of belonging to \mathcal{A}_ϕ . We will see that with the addition of a sufficient number of random clauses we can simulate the addition of x' .

Proof of Theorem 8. Suppose for a contradiction that \mathcal{A}_ϕ has a coarse threshold. Then there exist $\gamma, p_\gamma = o(1)$ and C as in Theorem 7, and so one of the two cases in its conclusion must hold.

Case 1: $\mu_p(x \in \{0, 1\}^n : x \text{ contains } x' \in \mathcal{A}_\phi \text{ of size } |x'| \leq 10C) > \delta$.

If the size of a formula x' is $\leq 10C$, then its clauses involve at most $10Ck$ variables. Since $x' \in \mathcal{A}_\phi$, and it is satisfied by v^0 , assigning the planted value to the variables appearing in x' and arbitrary values to the other variables yields a satisfying assignment. It follows that $Z(x') \geq 2^{n-10Ck} > 2^{\phi n}$ for large enough n , so $x' \notin \mathcal{A}_\phi$. This proves that Case 1 cannot occur.

Case 2: there exists $x' \notin \mathcal{A}_\phi$ of size $|x'| \leq 10C$ such that the conditional probability satisfies $\mu_{p_\gamma}(x \in \mathcal{A}_\phi | x' \subset x) > \gamma + \delta$.

Clearly x' is satisfied by v^0 . Denote by $t \leq 10Ck$ the number of variables appearing in x' . Without loss of generality, assume these variables are v_1, \dots, v_t . For a t -tuple $v = (v_{i_1}, \dots, v_{i_t})$ of distinct variables, we write $x'(v)$ to denote the result of relabeling each variable v_j in x' to v_{i_j} . Since \mathcal{A}_ϕ has permutation symmetry, it follows that for any t -tuple v , the conditional probability satisfies $\mu_p(x \in \mathcal{A}_\phi | x(v) \subset x) > \gamma + \delta$. We write x^* to mean the result of taking $x(v)$ after drawing a uniformly random t -tuple v . In other words, if a random formula $F^0(n, p_\gamma)$ is drawn, the union $F^0(n, p_\gamma) \cup x^*$ belongs to \mathcal{A}_ϕ with probability at least $\gamma + \delta$.

Now, since $p_\gamma \frac{d\mu_p(\mathcal{A}_\phi)}{dp} \big|_{p=p_\gamma} < C$ it follows that $\lim_{\varepsilon \rightarrow \infty} \frac{\mu_{p_\gamma + \varepsilon p_\gamma}(\mathcal{A}_\phi) - \mu_{p_\gamma}(\mathcal{A}_\phi)}{\varepsilon p_\gamma} < C$. Thus, for some ε we have $\mu_{p_\gamma + \varepsilon p_\gamma}(\mathcal{A}_\phi) < \gamma + \frac{\delta}{2}$. Further, (by a standard two-round exposure argument) choosing a formula $F^0(n, p_\gamma + \varepsilon p_\gamma)$ is equivalent to choosing formulae $F^0(n, p_\gamma)$ and $F^0(n, \varepsilon' p_\gamma)$ for some ε' and taking their union. Note that $\varepsilon, \varepsilon'$ don't depend on n , since C does not.

Denote by x^* a random copy of x' drawn as above. Then the above tells us that

$$\mathbb{P}[F^0(n, p_\gamma) \cup x^* \in \mathcal{A}_\phi] > \gamma + \delta$$

while

$$\mathbb{P}[F^0(n, p_\gamma) \cup F^0(n, \varepsilon' p_\gamma) \in \mathcal{A}_\phi] < \gamma + \frac{\delta}{2}.$$

It follows that for some formula $H_0 \in \{0, 1\}^N$ we have

$$\mathbb{P}[H_0 \cup x^* \in \mathcal{A}_\phi] - \mathbb{P}[H_0 \cup F^0(n, \varepsilon' p_\gamma) \in \mathcal{A}_\phi] > \frac{\delta}{2} \quad (4)$$

Clearly, $H_0 \notin \mathcal{A}_\phi$. Let's say that a t -tuple of distinct variables $v = (v_{i_1}, \dots, v_{i_t}) \in \{v_1, \dots, v_n\}^t$ is *bad* if $Z(H_0 \cup x(v)) < 2^{\phi n}$. It follows that at least a $\frac{\delta}{2}$ fraction of all $\binom{n}{t} t!$ t -tuples are bad. Let T be the set of bad tuples. We need the following theorem of Erdős and Simonovits [26].

Theorem 9 (Erdős and Simonovits). *Let k, t be positive integers and $0 \leq \gamma \leq 1$. There exists $\gamma' > 0$ such that for sufficiently large n , if $T \subset [n]^t$ is such that $|T| > \gamma n^t$ then with*

probability at least γ' a random choice of t disjoint k -tuples X_1, \dots, X_t from $[n]$ satisfies that every t -tuple (x_1, \dots, x_t) with $x_i \in X_i$ is bad. We say that X_1, \dots, X_t is T -complete.

We will obtain a contradiction from Theorem 9. Basically, we will ensure that with high probability, adding $F^0(n, \varepsilon'p)$ to H_0 implies adding clauses C_1, \dots, C_t , where each clause C_i forces some variable to be set to its planted value, and the set of k -tuples of variables in the clauses is T -complete.

Consider drawing t random clauses. Applying Theorem 9 with $\gamma = \frac{\delta}{2}$ we find some γ' for which the t k -clauses are T -complete with probability at least γ' . Given that they are T -complete, the probability that they are each of the form $\chi_e(v_{i_1} \dots v_{i_k}) \neq 1^k$ (in k -SAT or k -NAESAT case, or of the form $\chi_e(v_{i_1} \dots v_{i_k}) \neq k \pmod 2$ in the k -XORSAT case) is 2^{-kt} . Observe that each clause forces some variable to take the value 0, except in k -XORSAT when k is even and the clause forces some variable to take the value 1.

We claim that adding t such clauses to H_0 yields a formula with $< k^t 2^{\phi n}$ satisfying assignments. Indeed, suppose we have a satisfying assignment. Then at least one variable, say c_i , from each of the C_i must be set to 0 (1 in the even k -XORSAT case). This is at least as restrictive as containing $x((c_1, \dots, c_t))$, since $x(0^t)$ is satisfied (and in the even k -XORSAT case, therefore $x(1^t)$ is also satisfied). But (c_1, \dots, c_t) is a bad tuple so there are fewer than $2^{\phi n}$ ways to extend these to the remaining variables to get a satisfying assignment for H_0 .

With high probability, $F(\varepsilon'p_\gamma)$ has $\Theta(\varepsilon'p_\gamma \binom{n}{k} (2^k - 1)) \rightarrow \infty$ clauses. So if we draw $F^0(n, \varepsilon'p_\gamma)$ the probability that the clauses added don't include t clauses which force a 0 variable as above is at most about $(1 - \gamma' 2^{-kt})^{\varepsilon'p_\gamma \binom{n}{k} (2^k - 1)/t}$, which we can make as small as we like as $n \rightarrow \infty$. In particular, we can assume it is smaller than $\frac{\delta}{2}$. In the event that $F^0(n, \varepsilon'p_\gamma)$ does include these t clauses C_1, \dots, C_t , consider a satisfying assignment of $H_0 \cup C_1 \dots C_t$. The probability that it satisfies a randomly chosen k -clause is $(1 - 2^{-k})$. Therefore, in this case the expected value of $Z(H_0 \cup F^0(n, \varepsilon'p_\gamma))$ is at most $k^t 2^{\phi n} (1 - 2^{-k})^{|F^0(n, \varepsilon'p_\gamma)| - t} < 2^{\phi n}$ with high probability. Applying Markov's inequality, we can ensure that with probability greater than $1 - \frac{\delta}{2}$, the formula $H_0 \cup F^0(n, \varepsilon'p_\gamma) \in \mathcal{A}_\phi$, contradicting (4). This proves Case 2 cannot occur and completes the proof of the lemma. \square

5 Open problems

As mentioned in the introduction, it is not obvious that the conditions on the predicate χ used to obtain concentration (see Definition 1) are compatible with the conditions ruling out easily invertible functions [16]. The bottleneck here seems to be Hypothesis H, which at a high-level, translates the sub-additivity of the logarithm of the number of solutions (used to obtain concentration) into a local convexity property of the predicate χ . If the convexity property were in fact necessary, then this would be in conflict with the choice of predicates that seem to avoid the undesirable balanceness properties, making the problem curiously tensed between concentration and hardness. It would hence be interesting to show that this a limitation of our current proof technique, unless concentration has anything to do with hardness.

Another interesting question would be to obtain convergence rates for the convergence in probability. We obtain an exponential rate in Theorem 5 using martingale arguments, but this does not apply to our results relying on Bourgain.

Finally, the results in this paper are about the most basic properties of the solution space, namely, its cardinality. It would be interesting to understand rigorously finer properties of the solution space for planted models to deduce proper choices of the rate and predicates for the Goldreich one-way function.

Acknowledgement

We would like to thank R. Impagliazzo for suggesting the Goldreich one-way function model to the first author, as well as A. Montanari for stimulating discussions.

References

- [1] E. Abbe and A. Montanari. Conditional random fields, planted constraint satisfaction and entropy concentration. *To appear in the journal Theory of Computing, available at arXiv:1305.4274v2*. [2](#), [15](#)
- [2] E. Abbe and A. Montanari. On the concentration of the number of solutions of random satisfiability formulas. *Random Structures and Algorithms DOI 10.1002/rsa.20501, to appear*. arXiv:1006.3786v1, 2010. [1](#)
- [3] E. Abbe and A. Montanari. Conditional random fields, planted constraint satisfaction and entropy concentration. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 332–346. Springer, 2013. [1](#), [5](#)
- [4] E. Abbe and A. Montanari. On the concentration of the number of solutions of random satisfiability formulas. *Random Structures & Algorithms*, 2013. [1](#), [4](#), [5](#), [6](#), [14](#)
- [5] D. Achlioptas and A. Coja-Oghlan. Algorithmic barriers from phase transitions. In *Proceedings of the 2008 49th Annual IEEE Symposium on Foundations of Computer Science, FOCS '08*, pages 793–802, Washington, DC, USA, 2008. IEEE Computer Society. [1](#)
- [6] D. Achlioptas, A. Coja-Oghlan, and F. Ricci-Tersenghi. On the solution-space geometry of random constraint satisfaction problems. *Random Structures and Algorithms*, pages 251–268, 2010. [1](#), [8](#)
- [7] D. Achlioptas and E. Friedgut. A sharp threshold for k-colorability. *Random Structures and Algorithms*, 14(1):63–70, 1999. [8](#)
- [8] D. Achlioptas, H. Jia, and C. Moore. Hiding satisfying assignments: two are better than one. In *In Proceedings of AAAI04*, pages 131–136, 2004. [1](#)
- [9] D. Achlioptas, H. Kautz, and C. Gomes. Generating satisfiable problem instances. [1](#)
- [10] D. Achlioptas, J. H. Kim, M. Krivelevich, and P. Tetali. Two-coloring random hypergraphs. *Random Structures and Algorithms*, 20(2):249–259, 2002. [1](#)
- [11] D. Achlioptas, A. Naor, and Y. Peres. Rigorous Location of Phase Transitions in Hard Optimization Problems. *Nature*, 435:759–764, 2005. [1](#)
- [12] F. Altarelli, R. Monasson, and F. Zamponi. Can rare SAT formulas be easily recognized? On the efficiency of message passing algorithms for K-SAT at large clause-to-variable ratios. *Computing Research Repository*, abs/cs/060, 2006. [1](#)
- [13] K. Azuma. Weighted sums of certain dependent random variables. *Tohoku Mathematical Journal*, 19(3):357–367, 1967. [15](#)
- [14] W. Barthel, A. K. Hartmann, M. Leone, F. Ricci-Tersenghi, M. Weigt, and R. Zecchina. Hiding solutions in random satisfiability problems: A statistical mechanics approach. *Phys. Rev. Lett.*, 88:188701, Apr 2002. [1](#)

- [15] M. Bayati, D. Gamarnik, and P. Tetali. Combinatorial approach to the interpolation method and scaling limits in sparse random graphs. In *42nd Annual ACM Symposium on Theory of Computing*, pages 105–114, Cambridge, MA, June 2010. 1
- [16] A. Bogdanov and Y. Qiao. On the security of goldreich’s one-way function. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 392–405. Springer, 2009. 2, 6, 10
- [17] V. Chvátal and B. Reed. Mick gets some (the odds are on his side). *33th Annual Symposium on Foundations of Computer Science (Pittsburgh, PA, 1992)*, IEEE Comput. Soc. Press, pages 620–627, 1992. 1
- [18] A. Coja-Oghlan. The asymptotic k-SAT threshold. In *ACM Sympos. on Theory of Comput.*, New York, NY, 2014. 1
- [19] A. Coja-Oghlan, M. Krivelevich, and D. Vilenchik. Why almost all satisfiable k-cnf formulas are easy. In *Proceedings of the 13th International Conference on Analysis of Algorithms*, pages 89–102, 2007. 1
- [20] A. Coja-Oghlan and L. Zdeborová. The condensation transition in random hypergraph 2-coloring. arXiv:1107.2341, 2012. 2
- [21] H. Daudé and V. Ravelomanana. Random 2-xorsat at the satisfiability threshold. In *Proceedings of the 8th Latin American conference on Theoretical informatics, LATIN’08*, pages 12–23, Berlin, Heidelberg, 2008. Springer-Verlag. 1
- [22] W. F. de la Vega. On random 2-SAT. *manuscript*, 1992. 1
- [23] M. Dietzfelbinger, A. Goerdt, M. Mitzenmacher, A. Montanari, R. Pagh, and M. Rink. Tight thresholds for cuckoo hashing via XORSAT. Available at arXiv:0912.0287v1, 2010. 1
- [24] J. Ding, A. Sly, and N. Sun. Proof of the satisfiability conjecture for large k. arXiv:1411.0650, 2014. 1
- [25] O. Dubois and J. Mandler. The 3-xorsat threshold. In *Proceedings of the 43rd Symposium on Foundations of Computer Science*, FOCS ’02, pages 769–778, Washington, DC, USA, 2002. IEEE Computer Society. 1
- [26] P. Erdős and M. Simonovits. Supersaturated graphs and hypergraphs. *Combinatorica*, 3(2):181–192, 1983. 9
- [27] U. Feige, E. Mossel, and D. Vilenchik. Complete convergence of message passing algorithms for some satisfiability problems. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 339–350. Springer, 2006. 1
- [28] S. Franz and M. Leone. Replica bounds for optimization problems and diluted spin systems. *J. Stat. Phys.*, 111:535, 2003. 1
- [29] S. Franz, M. Leone, and F. Toninelli. Replica bounds for diluted non-Poissonian spin systems. *J. Phys. A*, 36:10967, 2003. 1
- [30] E. Friedgut. Hunting for sharp thresholds. *Random Structures & Algorithms*, 26(1-2):37–51, 2005. 7, 8

- [31] E. Friedgut and J. Bourgain. Sharp thresholds of graph properties, and the k -sat problem. *Journal of the American Mathematical Society*, 12(4):1017–1054, 1999. 1, 2, 7
- [32] E. Friedgut and G. Kalai. Every monotone graph property has a sharp threshold. *Proceedings of the American mathematical Society*, 124(10):2993–3002, 1996. 7
- [33] A. Goerdt. A Threshold for Unsatisfiability. *Journal of Computer and System Sciences*, 53:469–486, 1996. 1
- [34] O. Goldreich. Candidate one-way functions based on expander graphs. *IACR Cryptology ePrint Archive*, 2000:63, 2000. 4, 5
- [35] O. Goldreich. *Foundations of cryptography: a primer*, volume 1. Now Publishers Inc, 2005. 4
- [36] O. Goldreich. Studies in complexity and cryptography. chapter Candidate one-way functions based on expander graphs, pages 76–87. Springer-Verlag, Berlin, Heidelberg, 2011. 2
- [37] F. Guerra and F. L. Toninelli. The thermodynamic limit in mean field spin glasses. *Commun. Math. Phys.*, 230:71–79, 2002. 1
- [38] H. Haanpää, M. Järvisalo, P. Kaski, and I. Niemelä. Hard satisfiable clause sets for benchmarking equivalence reasoning techniques, 2005. 1
- [39] H. Jia, C. Moore, and D. Strain. Generating hard satisfiable formulas by hiding solutions deceptively. In *In AAAI*, pages 384–389. AAAI Press, 2005. 1
- [40] M. Krivelevich and A. Nachmias. Coloring complete bipartite graphs from random lists. *Random Structures & Algorithms*, 29(4):436–449, 2006. 7, 8
- [41] F. Krzakala, A. Montanari, F. Ricci-Tersenghi, G. Semerjian, and L. Zdeborová. Gibbs States and the Set of Solutions of Random Constraint Satisfaction Problems. *Proc. Natl. Acad. Sci.*, 104:10318–10323, 2007. 2
- [42] F. Krzakala and L. Zdeborová. Hiding quiet solutions in random constraint satisfaction problems. *Phys. Rev. Lett.*, 102:238701, Jun 2009. 1
- [43] A. Montanari. Personal communication, 2014. 2
- [44] A. Montanari, R. Restrepo, and P. Tetali. Reconstruction and Clustering in Random Constraint Satisfaction Problems. CoRR abs/0904.2751, 2009. 1
- [45] D. Panchenko and M. Talagrand. Bounds for diluted mean-field spin glass models. *Prob. Theor. Rel. Fields*, 130:319–336, 2004. 1
- [46] B. Pittel and G. B. Sorkin. The Satisfiability Threshold for k -XORSAT. *arXiv:1212.1905*, 2012. 1
- [47] L. Zdeborová and F. Krzakala. Phase transitions in the coloring of random graphs. *Phys. Rev. E*, 76:031131, Sep 2007. 2
- [48] L. Zdeborová and F. Krzakala. Quiet planting in the locked constraint satisfaction problems. *SIAM Journal on Discrete Mathematics*, 25(2):750–770, 2011. 1

A Freezing the threshold

In this section we prove Theorem 4. The proof essentially follows arguments in [4], but we give it here for completeness.

Proof of Theorem 4. For $\alpha \in [0, \alpha^*]$, let $\phi_s(\alpha)$ denote the limit of the sequence $\psi_n(\alpha) = \frac{1}{n} \mathbb{E}[\log Z(F(n, \alpha))]$ which converges almost surely by Theorem 2.

Let $\phi_0 = \phi_s(\alpha_0)$ for some α_0 . In view of Lemma 3 it is enough to show that the sequence $\alpha_n(\phi_0)$ obtained there converges (unless ϕ_0 takes one of countably many values). Suppose that it does not. Let

$$\underline{\alpha_0} = \liminf_{n \rightarrow \infty} \alpha_n(\phi_0)$$

and

$$\overline{\alpha_0} = \limsup_{n \rightarrow \infty} \alpha_n(\phi_0).$$

By assumption $\underline{\alpha_0}$ and $\overline{\alpha_0}$ disagree. Then we can choose increasing sequences $\{m_i\}_{i=1}^\infty$ and $\{n_i\}_{i=1}^\infty$ such that

$$\lim_{i \rightarrow \infty} \alpha_{m_i}(\phi_0) = \overline{\alpha_0}$$

and

$$\lim_{i \rightarrow \infty} \alpha_{n_i}(\phi_0) = \underline{\alpha_0}.$$

Let $\underline{\alpha_0} \leq \alpha \leq \overline{\alpha_0}$. Then for sufficiently large i there exists $\epsilon > 0$ such that

$$Q_{m_i}(\alpha, \phi_0) \leq Q_{m_i}(\alpha_{m_i}(\phi_0) - \epsilon, \phi_0) \rightarrow 0 \text{ as } i \rightarrow \infty$$

and

$$Q_{n_i}(\alpha, \phi_0) \geq Q_{n_i}(\alpha_{n_i}(\phi_0) + \epsilon, \phi_0) \rightarrow 1 \text{ as } i \rightarrow \infty.$$

Moreover since $\alpha \geq 0$ we have

$$Q_{m_i}(\alpha, \phi_0) = \mathbb{P} \left[Z(F(m_i, \alpha)) < 2^{m_i \phi_0} \right] = \mathbb{P} \left[\frac{1}{m_i} \log Z(F(m_i, \alpha)) < \phi_0 \right]$$

and so we have

$$\lim_{i \rightarrow \infty} \mathbb{E} \left[\frac{1}{m_i} \log Z(F(m_i, \alpha)) \right] \geq \phi_0$$

i.e.

$$\phi_s(\alpha) \geq \phi_0,$$

since the above expectation $\psi_n(\alpha)$ converges to $\phi_s(\alpha)$. A similar argument shows that

$$\lim_{i \rightarrow \infty} \mathbb{E} \left[\frac{1}{n_i} \log Z(F(n_i, \alpha)) \right] \leq \phi_0$$

and so

$$\phi_s(\alpha) \leq \phi_0.$$

It follows that the function ϕ_s is constant on $(\underline{\alpha_0}, \overline{\alpha_0})$. Since ϕ_s is non-increasing on $[0, \infty)$ it follows from Froda's theorem that there are countably many values ϕ_0 for which $\alpha_n(\phi_0)$ does not converge. This completes the proof. \square

We are now all set to prove our main theorem, which we restate now for convenience.

Theorem. *For every $k \geq 2$, there exist a countable set \mathcal{D} and a function $\phi_s : \mathbb{R}_{\geq 0} \rightarrow [0, 1]$*

such that for every $\alpha \notin \mathcal{D}$ and every $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} Q_n(\alpha, \phi_s(\alpha) - \epsilon) = 0$$

$$\lim_{n \rightarrow \infty} Q_n(\alpha, \phi_s(\alpha) + \epsilon) = 0$$

Proof of Theorem 1. Let ϕ_s be the function obtained in Lemma 3, and let \mathcal{D} denote the (countable) set of its discontinuities. Assume $\alpha \in [0, \alpha^*] \setminus \mathcal{D}$, and let $\phi_s(\alpha) = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}[\log Z(F(n, \alpha))]$ as in Theorem 2.

Lemma 3 implies that for some countable \mathcal{C} , the limit $A(\phi) = \lim_{n \rightarrow \infty} \alpha_n(\phi)$ exists for each $\phi \in \phi_s([0, \alpha^*]) \setminus \mathcal{C}$. Thus, there exists some $\epsilon' < \epsilon$ such that for $\phi^* := \phi_s(\alpha) - \epsilon'$ we have $\alpha_n(\phi^*)$ converges to a limit $A(\phi^*) > \alpha$. Therefore, there exists $\delta > 0$ such that

$$Q_n(\alpha, \phi_s(\alpha) - \epsilon) \leq Q_n(\alpha, \phi^*) \leq Q_n(\alpha_n(\phi^*) - \delta, \phi^*) \rightarrow 0.$$

It follows that $Q_n(\alpha, \phi_s(\alpha) - \epsilon) \rightarrow 0$ as $n \rightarrow \infty$. A symmetric argument shows that $Q_n(\alpha, \phi_s(\alpha) + \epsilon) \rightarrow 1$ as $n \rightarrow \infty$. This completes the proof. \square

B Proof of Theorem 5

Finally we now give the proof of Theorem 5.

Proof of Theorem 5. We consider the model of (3) for the Goldreich one-way function candidate. Let us denote by X a uniformly drawn input in $\{0, 1\}^n$ and by G a random hypergraph of fixed density. The output of the Goldreich one-way function candidate is the vector $Y(X, G) = \{\chi(X[e])\}_{e \in E(G)}$. We denote the number of pre-images of this output by $Z(X, G)$. In what follows, we show that the random variable $L(G) := \mathbb{E}_X \log Z(X, G)$ concentrates around its expectation $\mathbb{E}_G L(G)$, which depends on n .

For that purpose, we show that for any $\epsilon > 0$,

$$\mathbb{P}_G\{|L(G) - \mathbb{E}_G L(G)| \geq n\epsilon\} \leq 2e^{-n\epsilon^2/2}, \quad (5)$$

which implies that $|L(G)/n - \mathbb{E}_G L(G)/n|$ converges almost surely to 0 from the Borel-Cantelli Lemma. The above inequality results from a standard application of the Azuma-Hoeffding inequality [13], as used in [1] for more general models. Since the graph is Erdos-Renyi, we consider equivalently the edges to be drawn uniformly at random (conditioning on the number of edges in the graph). We need to show that, if e is an edge picked uniformly at random and $G \cup e$ is the augmented graph, the increment $L(G) - L(G \cup e)$ is bounded. In fact,

$$L(G) - L(G \cup e) = \mathbb{E}_X \log \frac{Z(X, G)}{Z(X, G \cup e)} \quad (6)$$

$$\leq \log \mathbb{E}_X \frac{Z(X, G)}{Z(X, G \cup e)} \quad (7)$$

$$= -\log \mathbb{E}_X \mathbb{E}_{U|X} \mathbb{1}(\chi(X[e]) = \chi(U[e])) \quad (8)$$

where U is a random vector uniformly drawn among all vectors $u \in \{0, 1\}^n$ such that the output of u is the same as the output of X on the one-way function defined by G and χ . Note that X and U are not independent but exchangeable, i.e., they are independent conditionally

on their common output Y . Therefore

$$\mathbb{E}_X \mathbb{E}_{U|X} \mathbf{1}(\chi(X[e]) = \chi(U[e])) = \mathbb{E}_{X,U} \mathbf{1}(\chi(X[e]) = \chi(U[e])) \quad (9)$$

$$= \mathbb{E}_{X[e], U[e]} \mathbf{1}(\chi(X[e]) = \chi(U[e])) \quad (10)$$

$$= \sum_y \mathbb{E}_{X[e], U[e]|Y=y} \mathbf{1}(\chi(X[e]) = \chi(U[e])) \mathbb{P}\{Y = y\} \quad (11)$$

$$= \sum_y (\mathbb{P}\{S_0|Y = y\}^2 + \mathbb{P}\{S_1|Y = y\}^2) \mathbb{P}\{Y = y\} \quad (12)$$

where $\mathbb{P}\{S_i|Y = y\}$ is the probability that $X[e]$ belongs to $\chi^{-1}(i)$ given that $Y = y$, for $i = 0, 1$. Since $\mathbb{P}\{S_0|Y = y\} + \mathbb{P}\{S_1|Y = y\} = 1$, we have $\mathbb{P}\{S_0|Y = y\}^2 + \mathbb{P}\{S_1|Y = y\}^2 \geq 1/2$, hence

$$\sum_y (\mathbb{P}\{S_0|Y = y\}^2 + \mathbb{P}\{S_1|Y = y\}^2) \mathbb{P}\{Y = y\} \geq 1/2, \quad (13)$$

and (8) is upper-bounded by $\log(2) = 1$. \square

C A sharp n and χ -dependent threshold for CSPs from Goldreich's functions

To prove Lemma 8 for a Goldreich random CSP as in (3) we need only make a slight modification to the proof in Section 4.2. First, to put ourselves in the setting where we have a product measure we fix the predicate χ (we do not need to fix the planted solution as we did above). This implies that the threshold we obtain may depend on both n and χ . As before, let $\mathcal{A}_\phi = \{F \in \{0, 1\}^N; Z(F) < 2^{\phi n}\}$, and now let $F = F(n, p)$ denote a CSP obtained as in (3), with χ fixed and denote by v^0 the planted solution. The space we are working in is $\{0, 1\}^N$, where $N = 2^{\binom{n}{k}}$, and indeed $\mu_p(F) = p^{|F|}(1-p)^{N-|F|}$.

Lemma 3 can be restated as follows.

Lemma 10. *For a fixed k and $\phi > 0$, the property \mathcal{A}_ϕ has a sharp threshold.*

The only place in which the proof of Lemma 10 differs from the proof of Lemma 8 is in the application of Theorem 9, but we give the details for completeness.

Proof of Theorem 10. Suppose for a contradiction that \mathcal{A}_ϕ has a coarse threshold. Then there exist $\gamma, p_\gamma = o(1)$ and C as in Theorem 7, and so one of the two cases in its conclusion must hold.

Case 1: $\mu_p(x \in \{0, 1\}^N : x \text{ contains } x' \in \mathcal{A}_\phi \text{ of size } |x'| \leq 10C) > \delta$.

If the size of a formula x' is $\leq 10C$, then its clauses involve at most $10Ck$ variables. Since $x' \in \mathcal{A}_\phi$, and it is satisfied by v^0 , assigning the planted value to the variables appearing in x' and arbitrary values to the other variables yields a satisfying assignment. It follows that $Z(x') \geq 2^{n-10Ck} > 2^{\phi n}$ for large enough n , so $x' \notin \mathcal{A}_\phi$. This proves that Case 1 cannot occur.

Case 2: there exists $x' \notin \mathcal{A}_\phi$ of size $|x'| \leq 10C$ such that the conditional probability satisfies $\mu_{p_\gamma}(x \in \mathcal{A}_\phi | x' \subset x) > \gamma + \delta$.

Clearly x' is satisfied by v^0 . Denote by $t \leq 10Ck$ the number of variables appearing in x' . Without loss of generality, assume these variables are v_1, \dots, v_t . For a t -tuple $v = (v_{i_1}, \dots, v_{i_t})$ of distinct variables, we write $x'(v)$ to denote the result of relabeling each variable v_j in x' to v_{i_j} . Since \mathcal{A}_ϕ has permutation symmetry, it follows that for any t -tuple v , the conditional probability satisfies $\mu_p(x \in \mathcal{A}_\phi | x(v) \subset x) > \gamma + \delta$. We write x^* to mean the result of taking $x(v)$

after drawing a uniformly random t -tuple v . In other words, if a random formula $F^0(n, p_\gamma)$ is drawn, the union $F^0(n, p_\gamma) \cup x^*$ belongs to \mathcal{A}_ϕ with probability at least $\gamma + \delta$.

Now, since $p_\gamma \frac{d\mu_p(\mathcal{A}_\phi)}{dp} \big|_{p=p_\gamma} < C$ it follows that $\lim_{\varepsilon \rightarrow \infty} \frac{\mu_{p_\gamma + \varepsilon p_\gamma}(\mathcal{A}_\phi) - \mu_{p_\gamma}(\mathcal{A}_\phi)}{\varepsilon p_\gamma} < C$. Thus, for some ε we have $\mu_{p_\gamma + \varepsilon p_\gamma}(\mathcal{A}_\phi) < \gamma + \frac{\delta}{2}$. Further, (by a standard two-round exposure argument) choosing a formula $F^0(n, p_\gamma + \varepsilon p_\gamma)$ is equivalent to choosing formulae $F^0(n, p_\gamma)$ and $F^0(n, \varepsilon' p_\gamma)$ for some ε' and taking their union. Note that $\varepsilon, \varepsilon'$ don't depend on n , since C does not.

Denote by x^* a random copy of x' drawn as above. Then the above tells us that

$$\mathbb{P}[F^0(n, p_\gamma) \cup x^* \in \mathcal{A}_\phi] > \gamma + \delta$$

while

$$\mathbb{P}[F^0(n, p_\gamma) \cup F^0(n, \varepsilon' p_\gamma) \in \mathcal{A}_\phi] < \gamma + \frac{\delta}{2}.$$

It follows that for some formula $H_0 \in \{0, 1\}^N$ we have

$$\mathbb{P}[H_0 \cup x^* \in \mathcal{A}_\phi] - \mathbb{P}[H_0 \cup F^0(n, \varepsilon' p_\gamma) \in \mathcal{A}_\phi] > \frac{\delta}{2} \quad (14)$$

Clearly, $H_0 \notin \mathcal{A}_\phi$. Let's say that a t -tuple of distinct variables $v = (v_{i_1}, \dots, v_{i_t}) \in \{v_1, \dots, v_n\}^t$ is *bad* if $Z(H_0 \cup x(v)) < 2^{\phi n}$. It follows that at least a $\frac{\delta}{2}$ fraction of all $\binom{n}{t} t!$ t -tuples are bad. Let T be the set of bad tuples. We need Erdős and Simonovits' Theorem 9 again.

We will ensure that with high probability, adding $F^0(n, \varepsilon' p)$ to H_0 implies adding clauses C_1, \dots, C_t , where each clause C_i forces some variable to be set to its planted value, and the set of k -tuples of variables in the clauses is T -complete.

Consider drawing t random clauses. Applying Theorem 9 with $\gamma = \frac{\delta}{2}$ we find some γ' for which the t k -clauses are T -complete with probability at least γ' . Given that they are T -complete, the probability that they are each of the form $\chi(v_{i_1} \dots v_{i_k}) = \chi(v_{i_1}^0 \dots v_{i_k}^0)$ By the antisymmetry of *chi*, each such clause forces some variable to take the planted value.

We claim that adding t such clauses to H_0 yields a formula with $< k^t 2^{\phi n}$ satisfying assignments. Indeed, suppose we have a satisfying assignment. Then at least one variable, say c_i , from each of the C_i must be set to the planted value c_i^0 . But (c_1, \dots, c_t) is a bad tuple so there are fewer than $2^{\phi n}$ ways to extend these to the remaining variables to get a satisfying assignment for H_0 .

With high probability, $F(\varepsilon' p_\gamma)$ has $\Theta(\varepsilon' p_\gamma \binom{n}{k} (2^k - 1)) \rightarrow \infty$ clauses. So if we draw $F^0(n, \varepsilon' p_\gamma)$ the probability that the clauses added don't include t clauses which force a 0 variable as above is at most about $(1 - \gamma' 2^{-kt})^{\varepsilon' p_\gamma \binom{n}{k} (2^k - 1)/t}$, which we can make as small as we like as $n \rightarrow \infty$. In particular, we can assume it is smaller than $\frac{\delta}{2}$. In the event that $F^0(n, \varepsilon' p_\gamma)$ does include these t clauses C_1, \dots, C_t , consider a satisfying assignment of $H_0 \cup C_1 \dots C_t$. The probability that it satisfies a randomly chosen k -clause is $(1 - 2^{-k})$. Therefore, in this case the expected value of $Z(H_0 \cup F^0(n, \varepsilon' p_\gamma))$ is at most $k^t 2^{\phi n} (1 - 2^{-k})^{|F^0(n, \varepsilon' p_\gamma)| - t} < 2^{\phi n}$ with high probability. Applying Markov's inequality, we can ensure that with probability greater than $1 - \frac{\delta}{2}$, the formula $H_0 \cup F^0(n, \varepsilon' p_\gamma) \in \mathcal{A}_\phi$, contradicting (14). This proves Case 2 cannot occur and completes the proof of the lemma. \square